

# PATENT ABSTRACTS OF JAPAN

(11)Publication number : 07-175698

(43)Date of publication of application : 14.07.1995

(51)Int.Cl.

G06F 12/00

(21)Application number : 05-317832

(71)Applicant : FUJITSU LTD

(22)Date of filing : 17.12.1993

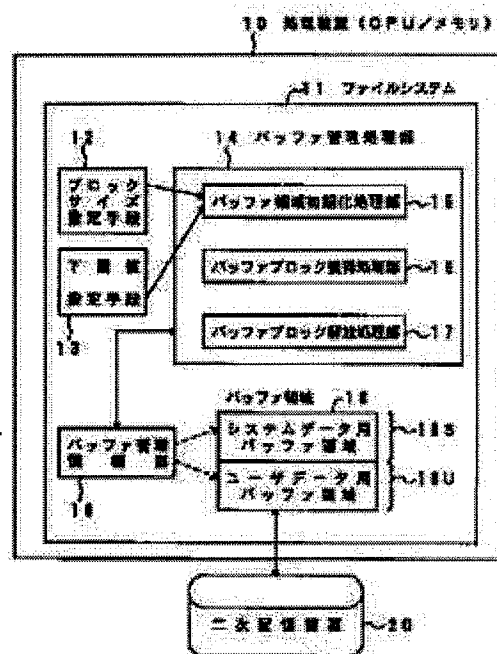
(72)Inventor : MURAKAMI TAKEO

## (54) FILE SYSTEM

### (57)Abstract:

**PURPOSE:** To improve input/output performance by dividing a buffer area for file input/output into areas by purposes including buffer areas for system data and user data, independently managing the respective areas, and controlling the possession and release of a requested buffer.

**CONSTITUTION:** Concerning the file system for a computer system provided with a central processing unit, main storage device and second-order storage device 20, a buffer area 18 for file input/output is initialized by a buffer area initialize processing part 15 and constituted so as to be divided into the plural areas by purposes such as a buffer area 18S for system data and a buffer area 18U for user data at this time of initializing. A buffer management processing part 14 independently manages the respective areas, and a buffer block possession processing part 16 and a buffer block release processing part 17 controls the possession and release of the requested buffer. Thus, input/output processing can be accelerated by utilizing the buffer for a parallel computer system or the like.



**\* NOTICES \***

JPO and INPIT are not responsible for any damages caused by the use of this translation.

- 1.This document has been translated by computer. So the translation may not reflect the original precisely.
- 2.\*\*\*\* shows the word which can not be translated.
- 3.In the drawings, any words are not translated.

---

**CLAIMS**

---

[Claim(s)]

[Claim 1]In a file system in a computer system provided with a central processing unit, a main memory unit, and a secondary memory, Divide the buffer space (18) into a buffer space for file input and output (18), and a field according to two or more uses which include a use of a buffer space for system data (18S), and a buffer space for user data (18U) at least, and each field independently. A file system provided with a buffer management treating part (14) which manages and controls acquisition and release of a buffer which were demanded.

[Claim 2]Have a block size setting means (12) which specifies size of a block which is an executive unit of a buffer in the file system according to claim 1 as every [ which was divided according to said use ] buffer space (18S, 18U), and Said buffer management treating part. A file system, wherein (14) is constituted so that acquisition and release of a buffer which managed a field and made the block size a unit with block size beforehand specified as said every divided buffer space (18S, 18U) may be controlled.

[Claim 3]In the file system according to claim 1, said buffer management treating part (14), A file system constituting so that it may change into a state which can assign a block in use and beyond a predetermined value may secure the number of assignable blocks, if the number of blocks which exist in a buffer space, and which can be assigned is less than a predetermined value.

[Claim 4]In the file system according to claim 1, a number of a block which can be assigned of lower limits which exist in each buffer space about a block which is an executive unit of a buffer to every [ which was divided according to said use ] buffer space (18S, 18U). Have a lower limit setting means (13) to specify, and said buffer management treating part (14), If the number of blocks which exist in a buffer space and which can be assigned is less than a block lower limit which was beforehand specified as said every divided buffer space (18S, 18U) and which can be assigned, A file system constituting so that it may change into a state which can assign a block in use and more than a block lower limit that can be assigned may secure the number of assignable blocks.

[Claim 5]In a file system in a computer system provided with a central processing unit, a main memory unit, and a secondary memory, A buffer block is acquired from a buffer space for file input and output (18), and said buffer space (18) to swap spaces. When it swaps out, set a dirty display which shows that the contents of the block have not been written in a secondary memory as information which manages the block, release the block to it, and to it for swapping in. A swap control means which cancels said dirty display and releases the block after swapping in when a buffer block is acquired, Said buffer space (18). By managing, and returning the contents of the block with which said dirty display is set up to a secondary memory, when the number of blocks which exist in a buffer space and which can be assigned is less than a block lower limit which was specified beforehand and which can be assigned, the number of assignable blocks. A file system provided with a buffer management treating part (14) which more than a block lower limit that can be assigned secures.

---

[Translation done.]

**\* NOTICES \***

JPO and INPIT are not responsible for any damages caused by the use of this translation.

- 1.This document has been translated by computer. So the translation may not reflect the original precisely.
- 2.\*\*\*\* shows the word which can not be translated.
- 3.In the drawings, any words are not translated.

---

**DETAILED DESCRIPTION**

---

[Detailed Description of the Invention]

[0001]

[Industrial Application]This invention relates to the file system which performs buffer management in consideration of input and output of system data, such as a swap.

[0002]In recent years, increase of the data volume treated in fields, such as large-scale numerical computation and time varying image processing, is remarkable, and the program itself has complicated and large-scale-sized it. Therefore, the computer system which can process more nearly mass data at high speed is called for. On the other hand, by new development of highly efficient CPU, maturation of parallel processing technology, etc., the improvement in computation capability is remarkable, and large scale-ization of a secondary storage and main memory is also following it. However, compared with these, the elongation of the data transfer rate of a secondary storage is small, and serves as a bottleneck of the whole system.

[0003]Then, main memory is used as a buffer of a secondary storage, and attaining improvement in the speed of radial transfer is usually performed. However, since there is a limit also in a main storage capacity, it is necessary to manage the limited buffer space efficiently.

[0004]

[Description of the Prior Art]In the conventional system, the whole buffer space was divided into the fixed-length block, and buffer management was performed. In a buffer space, the file which systems, such as a swapfile and dumping, use, and the file which users, such as a load module and a data file, use will be intermingled. Therefore, when the swap was performed or acquisition of dumping was performed, the buffer space might be occupied by the system, and a user program could not secure sufficient buffer space, but the problem that input-and-output performance fell might arise.

[0005]Since buffers with an enough system were not able to be secured if a swap is started when the job which outputs and inputs mass data runs and the buffer space is occupied, there was a problem that processing of a swap might become slow.

[0006]

[Problem(s) to be Solved by the Invention]As mentioned above, in the system which adopted the conventional buffer management system, Since one buffer space was used by the object for system data, and the object for user data, having competed, one of data occupied almost all buffer spaces, and there was a problem that the input-and-output performance of another side might fall extremely.

[0007]An object of this invention is to aim at improvement in input-and-output performance by dividing a buffer space according to a use and performing buffer management in consideration of the character of input and output.

[0008]

[Means for Solving the Problem]In the invention according to claim 1, as shown, for example in drawing 1, the buffer space 18 for file input and output is divided according to a use. A divided field

is called partition. One is the buffer space 18S for system data which systems secured when a job (process) is swapped out, such as a swapfile and dumping, use. Another is the buffer space 18U for user data which a user program uses. The buffer management treating part 14 manages each buffer space independently, and controls acquisition and release of a buffer which were demanded. The buffer space 18S for system data or the buffer space 18U for user data may be further divided according to a use.

[0009]When setting up a buffer at the time of starting of a system, it enables it to specify block size which is an executive unit of a buffer by the block size setting means 12 in the invention according to claim 2 as every [ which was divided according to a use ] buffer space (partition). The buffer management treating part 14 manages the buffer space 18S for system data, and the buffer space 18U for user data in block size different, respectively, and controls acquisition and release of a buffer which made the block size a unit by specified block size.

[0010]In the invention according to claim 3, the buffer management treating part 14 has managed the number of blocks which can be assigned for every partition. An assignable block is a block which is in agreement with the contents on a disk, although a block with which intact and effective data is not stored, or effective data is stored. Effective data is stored and the contents cannot assign a block (this is hereafter called dirty block) which is not reflected in a disk top.

[0011]When the buffer management treating part 14 performs block assignment, it investigates whether it is less than constant value (block lower limit which can be assigned) with the block count which can be assigned. If it has become less than constant value, the contents of the dirty block under list will be returned to the secondary memories 20, such as a disk, and the number of assignable blocks will be increased.

[0012]When setting up a buffer at the time of starting of a system, it enables it to specify a block lower limit which can be assigned by the lower limit setting means 13 in the invention according to claim 4 as every [ which was divided according to a use ] buffer space (partition). It is based on this specified block lower limit that can be assigned when the buffer management treating part 14 secures the block count in which assignment beyond constant value is possible.

[0013]In the invention according to claim 5, a context of an executing job swapped out or a process is written in a swapfile. At this time, if sufficient buffer space is secured to swapfiles, data will be written in a buffer, and a dirty flag which shows that a buffer block has not been written in to the secondary memory 20 is set. Even if a block with which a dirty flag was set is released, it assigns and it is impossible, but if required, it can be changed into a block which can be assigned by write return to the secondary memory 20 at any time. Then, if swapping in of a job is performed and a swapfile is read again, control which clears a dirty flag of an applicable buffer block will be performed.

[0014]

[Function]In the invention according to claim 1, since both buffer spaces 18S and 18U are divided even when the data input/output for systems and the data input/output for users are performed in parallel, it is guaranteed that all can secure a constant rate of buffer spaces. Therefore, each buffer spaces 18S and 18U for system data and for user data can be used, and it can output and input at a stable speed.

[0015]Since the character of input and output and the size of data differ between the data for systems, and the data for users greatly, they differ also in the executive unit of the optimal buffer for improving the utilization efficiency of a buffer, and input-and-output performance. In the invention according to claim 2, for every partition, specification of block size is possible and respectively suitable block size can be set up in consideration of the data of each partition, and input-output behavioral characteristics.

[0016]If a new buffer acquisition request comes when the assignable block count is lost in a buffer, after returning a dirty block to a disk, it will be necessary to assign, and the waiting time of buffer acquisition will increase. In the invention according to claim 3, since beyond constant value secures

the block which can always be assigned, there is no waiting time of buffer acquisition and input-and-output performance can be raised.

[0017] Since the size of input and output and an access pattern change with uses, the optimum values of the block lower limit which can be assigned differ for every partition. In the invention according to claim 4, since setting out of a different block lower limit for every partition which can be assigned can be performed, input-and-output performance can be raised by setting up a value suitable for the operating characteristic of each partition, respectively.

[0018] In swap control of an executing job or a process, once a swapfile is read by swapping in, it will not be read again. Therefore, when swapping in happens, even if the dirty block remains on the buffer, it is not necessary to return the contents to the secondary memories 20, such as a disk. In the invention according to claim 5, since a dirty flag is cleared using this character at the time of swapping in, write return on a useless disk can be prevented, and input-and-output performance can be raised. Since swapping out and swapping in are performed only on a buffer, without accessing to the secondary memory 20 when a buffer has a margin, high-speed swap control is possible.

[0019]

[Example] Drawing 1 shows the example of composition of this invention. 10 in drawing 1 A central processing unit (CPU). And a main memory unit. (Memory) etc. -- from -- the becoming processing unit and 11 the input and output to a file. The file system to perform, the block size setting means as which 12 specifies the block size of a buffer for every partition, the lower limit setting means which specifies a number of the block which can be assigned of lower limits with which 13 exists in a buffer space for every partition, and 14 a buffer. It is a buffer management treating part to manage, What consists of the buffer space initialization processing part 15, the buffer block acquisition treating part 16, the buffer block release treating part 17, etc., the buffer space where 18 was divided into the partition of the buffer space 18S for system data, and the buffer space 18U for user data, The buffer management information bureau where 19 holds the management information of the buffer space 18, and 20 express secondary memories, such as a disk unit.

[0020] The buffer space initialization processing part 15 is a processing means to initialize the buffer space 18 and to set the management information of the partition as the buffer management information bureau 19 according to the parameter specified for every partition by the block size setting means 12 and the lower limit setting means 13. The buffer block acquisition treating part 16 is a processing means to acquire a buffer block from an applicable partition to the acquisition request of the buffer from the demand origin which uses a buffer. The buffer block release treating part 17 is a processing means to release the block to the release request of the gained buffer.

[0021] As shown in drawing 1, the buffer space 18 for file input and output is initialized by the buffer space initialization processing part 15, and is constituted by the field according to use of plurality, such as the buffer space 18S for system data, and the buffer space 18U for user data, at the time of this initialization so that division is possible. The buffer management treating part 14 manages each field independently, and controls the acquisition and release of a buffer which were demanded by the buffer block acquisition treating part 16 and the buffer block release treating part 17.

[0022] Drawing 2 is an example explanatory view of composition of the file system in working example of this invention. In drawing 2, the thing of drawing 1 and a same sign to what is shown in drawing 1. Corresponding, program A-C, such as an application program with which 30A-30C perform the input and output to a file, the file management processing part which 31 receive the input/output request to a file and is processed, and 32 express the disk unit in which a file is stored.

[0023] The file system 11 comprises the buffer space 18 the file management processing part 31, the buffer management treating part 14, and for file input and output, etc. The file management processing part 31 is a treatment module which receives input/output request from an application program and performs a data area peculiar to a program, the buffer of the file system 11, and data input/output processing between the disk units 32.

[0024]The buffer management treating part 14 is a treatment module which manages the buffer space 18 used with the file system 11, and performs acquisition of a buffer, and release processing by a block unit by the demand from the file management processing part 31. The buffer management treating part 14 is divided into the buffer space initialization processing part 15, the buffer block acquisition treating part 16, and the buffer block release treating part 17.

[0025]According to a use, it divides and the buffer space 18 is managed, as shown in drawing 3 (A). In this example, the buffer space 18 is divided into three partitions, the partition 1 is assigned as a buffer space for system data, and the two partitions 2 and 3 are assigned as a buffer space for user data.

[0026]The block used as the executive unit of acquisition with each partition and release can specify the size beforehand for every partition. Since input and output with much data volume generally increase like a swap or dumping in the case of system data, in this example, large block size is set up as size of the buffer blocks B11-B13 in the partition 1 for system data. On the other hand, in the partition 2 for user data, since there are generally much input and output of a small unit, it is set up so that the size of the buffer blocks B21-B24 may become small.

[0027]The lower limit of the number of blocks which can be beforehand assigned at the time of initialization can be specified now to each partition. These information is managed for every partition as partition management information, as shown in drawing 3 (B). This partition management information is held in the buffer management information bureau 19 which shows drawing 1.

[0028]Drawing 4 is a management information explanatory view of the buffer block in working example of this invention. Hash cue and a free list are held for every partition as a data structure for managing the buffer space 18 other than partition management information to a block unit. 40 shown in drawing 4 by the block which can be assigned, and disk writing return. A hash table for the free list for managing the block which can be changed into the block which can be assigned, and 41 to manage hash cue, and 42 express the buffer block information which is the management information provided in each block correspondence.

[0029]As shown, for example in drawing 4 (A) for management of a buffer block, the hash table 41 which has the hash value 0 - a cue pointer corresponding to H for every partition is formed, and the buffer block information 42 into which effective data went for every hash value is held in the form of a list. A hash value is at this example. [File ID and the position in a file] are used as the key. File ID is file identification information given to a meaning within a system. The position in a file is the offset (displacement) information from the head in the file.

[0030]The free list 40 holds the buffer block information 42 used as the object seized when assigning a new buffer block in the form of a list. The buffer block information 42 established corresponding to each buffer block has the field where the information on the position in a file is set to file ID corresponding to the contents of the block, as shown, for example in drawing 4 (B). It has the field of the addresses in a buffer of the block (for example, offset or a block number from a partition head, etc.). It has a data valid flag which shows that the contents of the dirty flag which shows not having written the contents of the block in a disk to others, and its block are effective. It has a pointer for the chaining of hash cue, a pointer for the chaining of a free list, etc.

[0031]Next, the swap control using the buffer management system by this invention is explained. Drawing 5 is a swap control explanatory view in working example of this invention.

[0032]A swap is processing which vacates the storage area which evacuated to other fields and occupied till then the storage area which the job is using now for other jobs. In the case of this example, as shown in drawing 5 (A), the context 50 of a job is swapped out to the block B11 acquired to the buffer space 18, and the dirty flag of the block B11 is made one. And the buffer block information 42 of the block B11 is connected to the free list 40. At the time of swapping in, the block B11 is acquired by the acquisition request of a buffer block, and the contents are restored as the context 50 of a job. If swapping in is carried out, the dirty flag of the block B11 will be cleared.

[0033]When the number of the blocks of the buffer space 18 which can be assigned becomes less than a predetermined lower limit, the contents of the block B11 under swap are returned to the disk unit 32, and are treated as a block which can be assigned. At the time of swapping in, when the data of the block B11 is not effective, the new block B12 is acquired, the swapped context 50 is read from the disk unit 32, and the context 50 is restored to a new job field from the block B12.

[0034]The swapfile 51 will be realized by the buffer space 18 and the field of the disk unit 32 by making it above. Since a dirty flag is cleared when swapping in is performed, the write return to the useless disk unit 32 can be prevented. Namely, when there is sufficient margin for the buffer space 18, access to the disk unit 32 becomes unnecessary in the case of swapping out/swapping in.

[0035]The block with which the dirty flag is turned on is called dirty block. A dirty block is a block which is not reflecting the contents on a disk. In the free list 40 shown in drawing 5 (B), the buffer block information 42b and 42d and the block with --42x are dirty blocks. Other blocks are blocks which can be assigned, if the number of these blocks that can be assigned becomes smaller than the block lower limit which was set up for every partition and which can be assigned, the writing to the disk unit 32 of a dirty block will be performed, and a dirty flag will be cleared. Thereby, as shown in drawing 5 (C), the number of the blocks which can be assigned is managed so that it may become more than [ which can be assigned ] a block lower limit.

[0036]Drawing 6 is a figure showing the change state of the block under swap in working example of this invention. In the state which shows in drawing 6 (A), if a buffer block is acquired for swapping out, required number reservation is carried out, and after swapping out, the block for swapfiles will make a dirty flag one, and will be released. The state of the buffer block information 42 connected with the free list 40 comes to be shown in drawing 6 (B). If swapping in is performed, the dirty flag of the block used as a swapfile will be returned at OFF, and it will change in the state where it is shown in drawing 6 (C).

[0037]Drawing 7 is a buffer space initialization processing flow chart by working example of this invention. The buffer space initialization processing part 15 shown in drawing 2 is called from the file management processing part 31 at the time of initialization of a system. The input parameters of an argument are the information on a buffer space, block size, and a block lower limit that can be assigned (Step 70 of drawing 7). Step 71 divides the buffer space initialization processing part 15 first in the block size unit which had the object domain specified. Next, by Step 72, about each block, the buffer block information 42 of the contents shown in drawing 4 (B) is initialized, and it puts into the free list 40. In Step 73, the specified block lower limit which can be assigned is set up as partition management information. In Step 74, the number of whole blocks is initialized as the block count which can be assigned. In Step 75, partition ID is assigned and the partition ID is returned in the file management processing part 31 of a calling agency (Step 76).

[0038]By calling the buffer space initialization processing part 15 according to the use of a buffer space, the file management processing part 31 can divide into plurality the buffer space 18 used with the file system 11, and can set up a partition according to a use.

[0039]Drawing 8 is a buffer block acquisition processing flow chart by working example of this invention. The buffer block acquisition treating part 16 shown in drawing 2 is called from the file management processing part 31 at the time of the data input/output of each programs 30A-30B. Input parameters are partition ID which specifies the partition which is the acquisition target of a block, file ID for a data input/output, and the information on the position in a file (Step 800 of drawing 8).

[0040]The buffer block acquisition treating part 16 calculates the hash value which used specified file ID and the position in a file as the key by Step 801 first. Next, based on the calculated hash value, by Step 802, it asks for the entry of the hash table 41 in the applicable partition shown in drawing 4, and looks for hash cue.

[0041]If the applicable buffer block information 42 whose file ID and position in a file correspond is found with the judgment of Step 803, when it progresses to Step 810 and is not found, it progresses



to the following step 804.

[0042]When the applicable buffer block information 42 is not found, by Step 804, it searches for the free list 40 and the buffer block information 42 of the block which was found first and which can be assigned is removed from the free list 40. In Step 805, the buffer block information 42 of this block is removed from old hash cue, and it puts in to the new hash cue corresponding to the hash value which calculated file ID and the position in a file as a key. At the following step 806, the buffer block information 42, such as setting out of file ID, setting out of the position in a file, and OFF of a data valid flag, is initialized, and the block count which can be assigned is reduced at Step 807.

[0043]By Step 808, it judges whether the block count which can be assigned became smaller than the block lower limit which can be assigned, and if not small, it progresses to Step 813. If smaller than the block lower limit which can be assigned, by Step 809, it searches for the free list 40, and the contents of the block with which the dirty flag is turned on will be returned to a disk, the block will be changed into the state which can be assigned, and the block count which can be assigned will be increased. Then, it progresses to Step 813.

[0044]By the judgment of Step 803, when the applicable buffer block information 42 is found, it is judged by Step 810 whether the buffer block information 42 is among the free list 40. If it is among the free list 40, the buffer block information 42 will be removed from the free list 40 by Step 811.

[0045]Next, by Step 812, the data valid flag in the buffer block information 42 is made one, and it progresses to Step 813. In Step 813, it returns to the file management processing part 31 of a calling agency by making buffer block information 42 of the acquired block into a print-out.

[0046]Drawing 9 is a buffer block release processing flow chart by working example of this invention. The buffer block release treating part 17 shown in drawing 2 is called from the file management processing part 31, and performs release processing of the acquired buffer block. Input parameters are partition ID which specifies the partition with which the block to release exists, buffer block information, and a swapping in flag (Step 90 of drawing 9). Swapping in flags are one and a flag which is set up at OFF in the case of others, when releasing the buffer block with which swapping in was performed.

[0047]The buffer block release treating part 17 puts the buffer block information 42 of the block to release into the free list 40 of an applicable partition by Step 91 first. Next, by Step 92, a swapping in flag judges one or OFF, and when a swapping in flag is one (i.e., when it is block release at the time of swapping in), a dirty flag is cleared by Step 93.

[0048]In Step 94, ON and OFF of a dirty flag is judged, and if the dirty flag is off, the block count which can be assigned will be increased. Then, it returns to the file management processing part 31 of a calling agency at Step 96.

[0049]Drawing 10 is a swap processing flow chart by working example of this invention. Swap control of the job by the swapping in processing shown in the file management processing part 31 or the swapping out processing shown in drawing 10 (A) by the processing means of the swap control which carried out figures omitted abbreviated, and drawing 10 (B) which are shown in drawing 2 is performed.

[0050]In swapping out, first, by Step 101, partition ID of the buffer space 18S for system data, file ID of a swapfile, and the position in a file are specified, the buffer block acquisition treating part 16 is called, and the block for swapfiles is acquired. Next, the context 50 of the job swapped out to the acquired buffer block by Step 102 is written in, and the dirty flag in the buffer block information 42 of the block is made one by Step 103.

[0051]Next, a swapping in flag is cleared by Step 104, and the buffer block release treating part 17 is called by Step 105 by making partition ID, buffer block information, and a swapping in flag into an argument. By Step 106, if control returns from the buffer block release treating part 17, if it judges whether swapping out of all the contexts 50 was completed and there is an unsettled portion, it will return to Step 101 and processing will be similarly repeated about the following block. If it ends about all the contexts 50, processing of swapping out will be ended.

[0052]In swapping in, first, by Step 111, partition ID of the buffer space 18S for system data, file ID of a swapfile, and the position in a file for swapping in are specified, and the buffer block acquisition treating part 16 is called. If a buffer block is acquired, with reference to the buffer block information 42, ON and OFF of a data valid flag will be judged by Step 112. If a data valid flag is one, it will progress to Step 114.

[0053]If the data valid flag is off, the applicable context 50 will be read from the swapfile on a disk to the acquired block by Step 113. In Step 114, the contents of the buffer block are copied to a job field, and the context 50 is restored. Then, a swapping in flag is made one by Step 115, and the buffer block release treating part 17 is called by Step 116 by making partition ID, buffer block information, and a swapping in flag into an argument. By Step 117, if control returns from the buffer block release treating part 17, if it judges whether the swapping in of all the contexts 50 was completed and there is an unsettled portion, it will return to Step 111 and processing will be similarly repeated about the following block. If it ends about all the contexts 50, processing of swapping in will be ended.

[0054]

[Effect of the Invention]As explained above, according to this invention, it becomes possible to perform buffer management according to use which suited the use of each input and output, and largely contributes to the improved efficiency of the whole system. In the parallel computer system etc. which comprise a processing unit of a large number which it says especially are tens of sets – thousands of sets, for example, The art of managing efficiently the buffer space limited in each processing unit, and using it is needed, and according to this invention, it becomes possible to realize improvement in the speed of radial transfer which used the buffer in such a system, and effective use of a buffer space.

---

[Translation done.]

(19)日本国特許庁 (J P)

(12) 公 開 特 許 公 報 (A)

(11)特許出願公開番号

特開平7-175698

(43)公開日 平成7年(1995)7月14日

(51)Int.Cl.<sup>6</sup>  
G 0 6 F 12/00

識別記号 庁内整理番号  
5 1 4 A 8944-5B

F I

技術表示箇所

審査請求 未請求 請求項の数5 O L (全 10 頁)

(21)出願番号 特願平5-317832

(22)出願日 平成5年(1993)12月17日

(71)出願人 000005223

富士通株式会社

神奈川県川崎市中原区上小田中1015番地

(72)発明者 村上 岳生

神奈川県川崎市中原区上小田中1015番地

富士通株式会社内

(74)代理人 弁理士 小笠原 吉義 (外2名)

(54)【発明の名称】 ファイルシステム

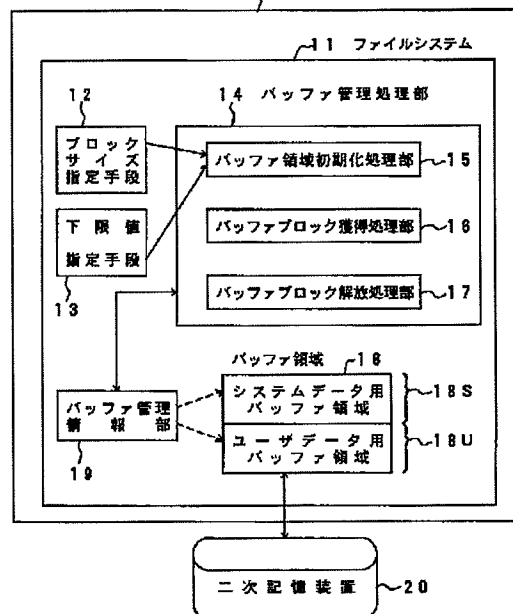
(57)【要約】

【目的】スワップ等のシステムデータの入出力を考慮してバッファ管理を行うファイルシステムに関し、入出力の各用途に適したバッファ管理により、入出力性能の向上を図ることを目的とする。

【構成】ファイル入出力用のバッファ領域18を、システムデータ用、ユーザデータ用というような用途別に複数の領域に分割し、バッファ管理処理部14により、それぞれの領域を独立に管理し、要求されたバッファの獲得および解放を制御する。また、用途別に分割したバッファ領域18S,18U 毎に、ブロックサイズおよび割り当て可能ブロック下限値を指定する手段を設け、領域の効率的利用とバッファ獲得の待ち時間の短縮を可能とする。

本 発 明 の 構 成 例

10 処理装置 (CPU/メモリ)



## 【特許請求の範囲】

【請求項1】 中央処理装置、主記憶装置および二次記憶装置を備えた計算機システムにおけるファイルシステムにおいて、ファイル入出力用のバッファ領域(18)と、そのバッファ領域(18)を、少なくともシステムデータ用バッファ領域(18S) およびユーザデータ用バッファ領域(18U) の用途を含む複数の用途別の領域に分割し、それぞれの領域を独立して管理し、要求されたバッファの獲得および解放を制御するバッファ管理処理部(14)とを備えたことを特徴とするファイルシステム。

【請求項2】 請求項1記載のファイルシステムにおいて、前記用途別に分割したバッファ領域(18S,18U) 毎に、バッファの管理単位であるブロックのサイズを指定するブロックサイズ指定手段(12)を備え、前記バッファ管理処理部(14)は、前記分割したバッファ領域(18S,18U) 毎にあらかじめ指定されたブロックサイズによって領域を管理し、そのブロックサイズを単位としたバッファの獲得および解放を制御するように構成されたことを特徴とするファイルシステム。

【請求項3】 請求項1記載のファイルシステムにおいて、前記バッファ管理処理部(14)は、バッファ領域内に存在する割り当て可能なブロックの数が所定値を下回ると、使用中のブロックを割り当て可能な状態に変更し、割り当て可能なブロックの数を所定値以上確保するように構成されたことを特徴とするファイルシステム。

【請求項4】 請求項1記載のファイルシステムにおいて、前記用途別に分割したバッファ領域(18S,18U) 毎に、バッファの管理単位であるブロックについて、それぞれのバッファ領域内に存在する割り当て可能なブロックの数の下限値を指定する下限値指定手段(13)を備え、前記バッファ管理処理部(14)は、バッファ領域内に存在する割り当て可能なブロックの数が前記分割したバッファ領域(18S,18U) 毎にあらかじめ指定された割り当て可能ブロック下限値を下回ると、使用中のブロックを割り当て可能な状態に変更し、割り当て可能なブロックの数を割り当て可能ブロック下限値以上確保するように構成されたことを特徴とするファイルシステム。

【請求項5】 中央処理装置、主記憶装置および二次記憶装置を備えた計算機システムにおけるファイルシステムにおいて、ファイル入出力用のバッファ領域(18)と、前記バッファ領域(18)からバッファブロックをスワップ領域用に獲得してスワップアウトをした際に、そのブロックを管理する情報に、ブロックの内容が二次記憶装置に未書き込みであることを示すダーティ表示を設定してそのブロックを解放し、スワップインのためにバッファブロックを獲得した際には、スワップイン後に前記ダーティ表示を解除してそのブロックを解放するスワップ制御手段と、前記バッファ領域(18)を管理し、バッファ領域内に存在する割り当て可能なブロックの数があらかじめ指定された割り当て可能ブロック下限値を下回った場

合に、前記ダーティ表示が設定されているブロックの内容を二次記憶装置に書き戻すことにより、割り当て可能なブロックの数を割り当て可能ブロック下限値以上確保するバッファ管理処理部(14)とを備えたことを特徴とするファイルシステム。

## 【発明の詳細な説明】

## 【0001】

【産業上の利用分野】本発明は、スワップ等のシステムデータの入出力を考慮してバッファ管理を行うファイルシステムに関する。

【0002】近年、大規模数値計算や動画処理等の分野において扱うデータ量の増大は著しく、またプログラム自体も複雑化し大容量化している。従って、より大容量のデータを高速に処理することが可能な計算機システムが求められている。一方、高性能なCPUの新規開発、並列処理技術の成熟等により、計算処理能力の向上は著しく、また二次記憶、主記憶の大容量化も進んでいる。しかし、これらに比べて二次記憶のデータ転送速度の伸びは小さく、システム全体のボトルネックとなっている。

【0003】そこで通常、主記憶を二次記憶のバッファとして利用し、入出力処理の高速化を図ることが行われている。しかし、主記憶容量にも限界があるので、限られたバッファ領域を効率良く管理する必要がある。

## 【0004】

【従来の技術】従来のシステムにおいては、バッファ領域全体を固定長のブロックに分割してバッファ管理を行っていた。バッファ領域はスワップファイル、ダンプ等のシステムが使用するファイルと、ロードモジュール、データファイル等のユーザが使用するファイルとが混在することになる。従って、スワップが行われたり、ダンプの取得が行われたりすると、バッファ領域がシステムに占有され、ユーザプログラムが十分なバッファ領域を確保できず、入出力性能が低下するという問題が生じることがあった。

【0005】また、大容量のデータの入出力を行うジョブが走行し、バッファ領域を占有している際に、スワップを起動すると、システムが十分なバッファを確保できないため、スワップの処理が遅くなることがあるという問題があった。

## 【0006】

【発明が解決しようとする課題】以上のように、従来のバッファ管理方式を採用したシステムにおいては、1つのバッファ領域をシステムデータ用とユーザデータ用とで競合して使用しているので、いずれかのデータがほとんどのバッファ領域を占有してしまい、他方の入出力性能が極端に落ちてしまうことがあるという問題があった。

【0007】本発明は、用途別にバッファ領域を分割し、入出力の性質を考慮したバッファ管理を行うことに

10

20

30

40

50

より、入出力性能の向上を図ることを目的とする。

#### 【0008】

【課題を解決するための手段】請求項1記載の発明では、例えば図1に示すように、ファイル入出力用のバッファ領域18を用途別に分割する。分割された領域をパーティションという。1つはジョブ（プロセス）がスワップアウトされる際に確保するスワップファイルやダンプ等のシステムが使用するシステムデータ用バッファ領域18Sである。もう1つはユーザプログラムが使用するユーザデータ用バッファ領域18Uである。バッファ管理処理部14は、それぞれのバッファ領域を独立して管理し、要求されたバッファの獲得および解放を制御する。なお、システムデータ用バッファ領域18Sまたはユーザデータ用バッファ領域18Uをさらに用途別に分割してもよい。

【0009】請求項2記載の発明では、システムの起動時にバッファの設定を行う際に、ブロックサイズ指定手段12によって、用途別に分割したバッファ領域（パーティション）毎に、バッファの管理単位であるブロックサイズを指定できるようにする。バッファ管理処理部14は、指定されたブロックサイズによってシステムデータ用バッファ領域18Sおよびユーザデータ用バッファ領域18Uを、それぞれ異なるブロックサイズで管理し、そのブロックサイズを単位としたバッファの獲得および解放を制御する。

【0010】請求項3記載の発明では、バッファ管理処理部14は、パーティション毎に割り当て可能なブロックの数を管理している。割り当て可能なブロックとは、未使用で有効なデータが格納されていないブロック、あるいは有効なデータが格納されているがディスク上の内容と一致しているブロックである。有効なデータが格納されており、その内容がディスク上に未反映のブロック（以下、これをダーティブロックという）は割り当て不可能である。

【0011】バッファ管理処理部14は、ブロック割り当てを行う際に、割り当て可能ブロック数がある一定値（割り当て可能ブロック下限値）より少ないかどうかを調べる。一定値より少なくなっていればリスト中のダーティブロックの内容をディスク等の二次記憶装置20に書き戻し、割り当て可能なブロックの数を増やす。

【0012】請求項4記載の発明では、システムの起動時にバッファの設定を行う際に、下限値指定手段13により、用途別に分割したバッファ領域（パーティション）毎に、割り当て可能ブロック下限値を指定できるようにする。バッファ管理処理部14が一定値以上の割り当て可能なブロック数を確保する際には、この指定された割り当て可能ブロック下限値を基準とする。

【0013】請求項5記載の発明では、スワップアウトされる実行中ジョブまたはプロセスのコンテキストがスワップファイルに書き込まれる。このとき、十分なバッ

ファ領域がスワップファイル用に確保されればデータをバッファに書き込み、バッファブロックが二次記憶装置20に対して未書き込みであることを示すダーティフラグをセットする。ダーティフラグがセットされたブロックは解放されても割り当て不可能であるが、必要であれば二次記憶装置20への書き戻しにより、いつでも割り当て可能ブロックに変更することができる。その後、ジョブのスワップインが行われ再びスワップファイルが読み込まれると、該当バッファブロックのダーティフラグをオフにする制御を行う。

#### 【0014】

【作用】請求項1記載の発明では、システム用のデータ入出力とユーザ用のデータ入出力とが並行して行われる場合でも、両者のバッファ領域18S、18Uが分割されているため、いずれも一定量のバッファ領域を確保できることが保証される。従って、システムデータ用、ユーザデータ用の各バッファ領域18S、18Uを使用して、安定した速度で入出力を行うことができる。

【0015】システム用のデータとユーザ用のデータとは、入出力の性質やデータのサイズが大きく異なるので、バッファの利用効率、入出力性能を上げるための最適なバッファの管理単位も異なる。請求項2記載の発明では、パーティション毎にブロックサイズの指定が可能であり、各パーティションのデータ、入出力特性を考慮してそれぞれ適切なブロックサイズを設定することができる。

【0016】割り当て可能なブロック数がバッファ中になくなった際に、新たなバッファ獲得要求が来ると、ダーティなブロックをディスクに書き戻してから割り当てる必要があり、バッファ獲得の待ち時間が増えてしまう。請求項3記載の発明では、常に割り当て可能なブロックを一定値以上確保しておくので、バッファ獲得の待ち時間がなく、入出力性能を向上させることができる。

【0017】また、用途によって入出力のサイズ、アクセスパターンは異なるため、パーティション毎に割り当て可能ブロック下限値の最適値は異なる。請求項4記載の発明では、パーティション毎に異なる割り当て可能ブロック下限値の設定ができるため、各パーティションの使用特性に適した値をそれぞれ設定することにより、入出力性能を向上させることができる。

【0018】実行中ジョブまたはプロセスのスワップ制御において、スワップファイルが一度スワップインで読み込まれると、再び読み込まれることはない。従って、スワップインが起こった時点でダーティなブロックがバッファ上に残っていても、その内容をディスク等の二次記憶装置20へ書き戻す必要はない。請求項5記載の発明では、この性質を利用してスワップイン時にダーティフラグをオフにするので無駄なディスクへの書き戻しを防ぐことができ、入出力性能を向上させることができる。また、バッファに余裕がある場合には、二次記憶装

置20へアクセスすることなく、バッファ上だけでスワップアウト、スワップインが行われるので、高速なスワップ制御が可能である。

#### 【0019】

【実施例】図1は本発明の構成例を示す。図1において、10は中央処理装置（CPU）および主記憶装置（メモリ）等からなる処理装置、11はファイルへの入出力を行うファイルシステム、12はパーティション毎にバッファのブロックサイズを指定するブロックサイズ指定手段、13はパーティション毎にバッファ領域内に存在する割り当て可能なブロックの数の下限値を指定する下限値指定手段、14はバッファを管理するバッファ管理処理部であって、バッファ領域初期化処理部15、バッファブロック獲得処理部16、バッファブロック解放処理部17などからなるもの、18はシステムデータ用バッファ領域18Sとユーザデータ用バッファ領域18Uのパーティションに分割されたバッファ領域、19はバッファ領域18の管理情報を保持するバッファ管理情報部、20はディスク装置等の二次記憶装置を表す。

【0020】バッファ領域初期化処理部15は、パーティション毎にブロックサイズ指定手段12および下限値指定手段13で指定されたパラメータに従って、バッファ領域18の初期化を行い、そのパーティションの管理情報をバッファ管理情報部19に設定する処理手段である。バッファブロック獲得処理部16は、バッファを使用する要求元からのバッファの獲得要求に対して、該当するパーティションからバッファブロックを獲得する処理手段である。バッファブロック解放処理部17は、獲得されたバッファの解放要求に対して、そのブロックを解放する処理手段である。

【0021】図1に示すように、ファイル入出力用のバッファ領域18は、バッファ領域初期化処理部15によって初期化され、この初期化時にシステムデータ用バッファ領域18Sおよびユーザデータ用バッファ領域18U等の複数の用途別の領域に分割可能に構成されている。バッファ管理処理部14は、それぞれの領域を独立して管理し、バッファブロック獲得処理部16およびバッファブロック解放処理部17によって、要求されたバッファの獲得および解放を制御する。

【0022】図2は、本発明の実施例におけるファイルシステムの構成例説明図である。図2において、図1と同符号のものは図1に示すものに対応し、30A~30Cはファイルへの入出力を行うアプリケーションプログラム等のプログラムA~C、31はファイルへの入出力要求を受け取って処理するファイル管理処理部、32はファイルが格納されるディスク装置を表す。

【0023】ファイルシステム11は、ファイル管理処理部31とバッファ管理処理部14とファイル入出力用のバッファ領域18などから構成されている。ファイル管理処理部31は、アプリケーションプログラムから入

出力要求を受け取り、プログラム固有のデータ領域、ファイルシステム11のバッファ、ディスク装置32の間のデータ入出力処理を行う処理モジュールである。

【0024】バッファ管理処理部14は、ファイルシステム11で使用するバッファ領域18を管理し、ファイル管理処理部31からの要求により、ブロック単位でバッファの獲得、解放処理を行う処理モジュールである。バッファ管理処理部14は、バッファ領域初期化処理部15と、バッファブロック獲得処理部16と、バッファブロック解放処理部17とに分かれている。

【0025】バッファ領域18は、図3（A）に示すように、用途別に分割して管理される。この例では、バッファ領域18は3つのパーティションに分割され、システムデータ用のバッファ領域としてパーティション1が、ユーザデータ用のバッファ領域として2つのパーティション2、3が割り当てられている。

【0026】各パーティションでの獲得、解放の管理単位となるブロックは、パーティション毎にあらかじめそのサイズを指定することができる。システムデータの場合、スワップやダンプ等のように一般にデータ量の多い入出力が多くなることから、この例では、システムデータ用のパーティション1におけるバッファブロックB11~B13のサイズとして、大きいブロックサイズが設定されている。一方、ユーザデータ用のパーティション2では、一般に小さい単位の入出力が多いことから、バッファブロックB21~B24のサイズが小さくなるように設定されている。

【0027】また、各パーティションに対して、初期化時にあらかじめ割り当て可能なブロックの数の下限値を指定できるようになっている。これらの情報は、図3（B）に示すように、パーティション管理情報として各パーティション毎に管理される。このパーティション管理情報は、図1に示すバッファ管理情報部19に保持される。

【0028】図4は本発明の実施例におけるバッファブロックの管理情報説明図である。パーティション管理情報の他に、バッファ領域18をブロック単位に管理するためのデータ構造として、各パーティション毎にハッシュキュー、フリーリストを保持している。図4に示す40は割り当て可能ブロックおよびディスク書き戻しによって割り当て可能ブロックに変更できるブロックを管理するためのフリーリスト、41はハッシュキューを管理するためのハッシュテーブル、42は各ブロック対応に設けられる管理情報であるバッファブロック情報を表す。

【0029】バッファブロックの管理のために、例えば図4（A）に示すように、各パーティション毎にハッシュ値0~Hに対応したキューポインタを持つハッシュテーブル41が設けられ、ハッシュ値毎に有効なデータの入ったバッファブロック情報42をリストの形で保持す

る。ハッシュ値は、この例では〔ファイルID、ファイル内位置〕をキーとしている。ファイルIDはシステム内で一意に付与されるファイル識別情報である。ファイル内位置は、そのファイルにおける先頭からのオフセット（変位）情報である。

【0030】フリーリスト40は、新規のバッファブロックを割り当てる際に横取りする対象となるバッファブロック情報42をリストの形で保持している。個々のバッファブロックに対応して設けられるバッファブロック情報42は、例えば図4（B）に示すように、そのブロックの内容に対応するファイルIDと、ファイル内位置の情報が設定されるフィールドを持つ。また、そのブロックのバッファ内アドレス（例えば、パーティション先頭からのオフセットまたはブロック番号等）のフィールドを持つ。他に、ブロックの内容がディスクに未書き込みであることを示すダーティフラグ、そのブロックの内容が有効であることを示すデータ有効フラグを持つ。さらに、ハッシュキューのチェイニングのためのポインタおよびフリーリストのチェイニングのためのポインタ等を持つ。

【0031】次に、本発明によるバッファ管理方式を用いたスワップ制御について説明する。図5は、本発明の実施例におけるスワップ制御説明図である。

【0032】スワップは、ジョブが現在使用している記憶領域を他の領域に退避して、それまで占めていた記憶領域を他のジョブ等のために明け渡す処理である。本実施例の場合、図5（A）に示すように、ジョブのコンテキスト50をバッファ領域18に獲得したブロックB11へスワップアウトし、ブロックB11のダーティフラグをオンにする。そして、そのブロックB11のバッファブロック情報42をフリーリスト40に接続しておく。スワップインのときには、バッファブロックの獲得要求によりブロックB11を獲得し、その内容をジョブのコンテキスト50として復元する。スワップインをしたならば、そのブロックB11のダーティフラグをオフにする。

【0033】スワップ中のブロックB11の内容は、バッファ領域18の割り当て可能ブロックの数が所定の下限值より少なくなった場合には、ディスク装置32に書き戻され、割り当て可能ブロックとして扱われる。スワップインのとき、ブロックB11のデータが有効でない場合には、新しいブロックB12を獲得して、スワップされたコンテキスト50をディスク装置32から読み出し、そのブロックB12から新たなジョブ領域にコンテキスト50を復元する。

【0034】以上のようにすることにより、スワップファイル51は、バッファ領域18とディスク装置32の領域とによって実現されることになる。スワップインが行われたときには、ダーティフラグをオフにするので、無駄なディスク装置32への書き戻しを防ぐことができ

る。すなわち、バッファ領域18に十分な余裕がある場合には、スワップアウト／スワップインの際にディスク装置32へのアクセスは不要となる。

【0035】ダーティフラグがオンになっているブロックをダーティブロックという。ダーティブロックは、ディスク上に内容を反映していないブロックである。図5（B）に示すフリーリスト40では、バッファブロック情報42b、42d、…42xを持つブロックがダーティブロックである。他のブロックは、割り当て可能ブロックであり、この割り当て可能ブロックの数が、パーティション毎に設定された割り当て可能ブロック下限値より小さくなると、ダーティブロックのディスク装置32への書き込みが行われ、ダーティフラグがオフにされる。これにより、図5（C）に示すように、割り当て可能ブロックの数が、割り当て可能ブロック下限値以上になるように管理される。

【0036】図6は、本発明の実施例におけるスワップ中のブロックの状態遷移を示す図である。図6（A）に示す状態で、スワップアウトのためにバッファブロックが獲得されると、スワップファイル用ブロックが必要個数確保され、スワップアウトの後、ダーティフラグをオンにして解放される。フリーリスト40につながるバッファブロック情報42の状態は、図6（B）に示すようになる。スワップインが行われると、スワップファイルとして用いられたブロックのダーティフラグがオフに戻され、図6（C）に示すような状態に遷移する。

【0037】図7は、本発明の実施例によるバッファ領域初期化処理フローチャートである。図2に示すバッファ領域初期化処理部15は、システムの初期化時にファイル管理処理部31から呼び出される。引数の入力パラメータは、バッファ領域の情報、ブロックサイズ、割り当て可能ブロック下限値である（図7のステップ70）。バッファ領域初期化処理部15は、まずステップ71により、対象領域を指定されたブロックサイズ単位で分割する。次にステップ72により、各ブロックについて、図4（B）に示す内容のバッファブロック情報42を初期化し、フリーリスト40に入れる。ステップ73では、指定された割り当て可能ブロック下限値を、パーティション管理情報として設定する。ステップ74では、全ブロック数を割り当て可能ブロック数として初期化する。ステップ75では、パーティションIDを割り当て、そのパーティションIDを呼び出し元のファイル管理処理部31に返却する（ステップ76）。

【0038】ファイル管理処理部31は、バッファ領域の用途別にバッファ領域初期化処理部15を呼び出すことにより、ファイルシステム11で使用するバッファ領域18を複数個に分割し、用途別にパーティションを設定することができる。

【0039】図8は、本発明の実施例によるバッファブロック獲得処理フローチャートである。図2に示すバ

バッファブロック獲得処理部16は、各プログラム30A～30Bのデータ入出力時にファイル管理処理部31から呼び出される。入力パラメータは、ブロックの獲得対象となるパーティションを指定するパーティションID、データ入出力対象のファイルID、ファイル内位置の情報である(図8のステップ800)。

【0040】バッファブロック獲得処理部16は、まずステップ801により、指定されたファイルIDとファイル内位置をキーとしたハッシュ値を計算する。次に、計算したハッシュ値をもとに、ステップ802により、図4に示す該当パーティションにおけるハッシュテーブル41のエントリを求め、ハッシュキューを探索する。

【0041】ステップ803の判定により、ファイルIDおよびファイル内位置が一致する該当バッファブロック情報42が見つかったならば、ステップ810へ進み、見つからなかった場合には、次のステップ804へ進む。

【0042】該当バッファブロック情報42が見つからなかった場合、ステップ804により、フリーリスト40を探索し、最初に見つかった割り当て可能ブロックのバッファブロック情報42をフリーリスト40から外す。ステップ805では、古いハッシュキューからこのブロックのバッファブロック情報42を外し、ファイルIDおよびファイル内位置をキーとして計算したハッシュ値に対応する新しいハッシュキューへ入れる。次のステップ806で、ファイルIDの設定、ファイル内位置の設定、データ有効フラグのオフなどのバッファブロック情報42の初期化を行い、ステップ807で割り当て可能ブロック数を減らす。

【0043】ステップ808により、割り当て可能ブロック数が割り当て可能ブロック下限値より小さくなったかどうかを判定し、小さくなければステップ813へ進む。割り当て可能ブロック下限値より小さくければ、ステップ809により、フリーリスト40を探索し、ダーティフラグがオンになっているブロックの内容をディスクへ書き戻し、そのブロックを割り当て可能状態にして、割り当て可能ブロック数を増やす。その後、ステップ813へ進む。

【0044】ステップ803の判定で、該当バッファブロック情報42が見つかった場合、ステップ810により、そのバッファブロック情報42がフリーリスト40中かどうかを判定する。フリーリスト40中であれば、ステップ811によりそのバッファブロック情報42をフリーリスト40から外す。

【0045】次に、ステップ812により、そのバッファブロック情報42におけるデータ有効フラグをオンにして、ステップ813へ進む。ステップ813では、獲得したブロックのバッファブロック情報42を出力情報として、呼び出し元のファイル管理処理部31へ復帰する。

【0046】図9は、本発明の実施例によるバッファブロック解放処理フローチャートである。図2に示すバッファブロック解放処理部17は、ファイル管理処理部31から呼び出され、獲得したバッファブロックの解放処理を行う。入力パラメータは、解放するブロックの存在するパーティションを指定するパーティションID、バッファブロック情報、スワップインフラグである(図9のステップ90)。スワップインフラグは、スワップインが行われたバッファブロックを解放する場合にオン、その他の場合にはオフに設定されるフラグである。

【0047】バッファブロック解放処理部17は、まずステップ91により、該当パーティションのフリーリスト40に、解放するブロックのバッファブロック情報42を入れる。次にステップ92により、スワップインフラグがオンかオフかを判定し、スワップインフラグがオンの場合、すなわちスワップイン時のブロック解放の場合には、ステップ93によってダーティフラグをオフにする。

【0048】ステップ94では、ダーティフラグのオン/オフを判定し、ダーティフラグがオフであれば割り当て可能ブロック数を増やす。その後、ステップ96で呼び出し元のファイル管理処理部31へ復帰する。

【0049】図10は、本発明の実施例によるスワップ処理フローチャートである。図2に示すファイル管理処理部31または図示省略したスワップ制御の処理手段によって、図10(A)に示すスワップアウト処理および図10(B)に示すスワップイン処理によるジョブのスワップ制御が行われる。

【0050】スワップアウトでは、まずステップ101により、システムデータ用バッファ領域18SのパーティションIDと、スワップファイルのファイルIDと、ファイル内位置とを指定して、バッファブロック獲得処理部16を呼び出し、スワップファイル用ブロックを獲得する。次に、ステップ102により、獲得したバッファブロックにスワップアウトするジョブのコンテキスト50を書き込み、ステップ103によって、そのブロックのバッファブロック情報42におけるダーティフラグをオンにする。

【0051】次に、ステップ104により、スワップインフラグをオフにし、ステップ105により、パーティションIDとバッファブロック情報とスワップインフラグとを引数として、バッファブロック解放処理部17を呼び出す。バッファブロック解放処理部17から制御が戻ったならば、ステップ106により、全コンテキスト50のスワップアウトが終了したかどうかを判定し、未処理部分があれば、ステップ101へ戻って、次のブロックについて同様に処理を繰り返す。全コンテキスト50について終了したならば、スワップアウトの処理を終了する。

【0052】スワップインでは、まずステップ111に



より、システムデータ用バッファ領域18SのパーティションIDと、スワップファイルのファイルIDと、スワップイン対象のファイル内位置とを指定して、バッファブロック獲得処理部16を呼び出す。バッファブロックを獲得したならば、ステップ112により、そのバッファブロック情報42を参照し、データ有効フラグのオン/オフを判定する。データ有効フラグがオンであれば、ステップ114へ進む。

【0053】データ有効フラグがオフであれば、ステップ113により、獲得したブロックにディスク上のスワップファイルから該当コンテキスト50を読み出す。ステップ114では、バッファブロックの内容をジョブ領域へ複写し、コンテキスト50の復元を行う。その後、ステップ115により、スワップインフラグをオンにし、ステップ116により、パーティションIDとバッファブロック情報とスワップインフラグとを引数として、バッファブロック解放処理部17を呼び出す。バッファブロック解放処理部17から制御が戻ったならば、ステップ117により、全コンテキスト50のスワップインが終了したかどうかを判定し、未処理部分があれば、ステップ111へ戻って、次のブロックについて同様に処理を繰り返す。全コンテキスト50について終了したならば、スワップインの処理を終了する。

#### 【0054】

【発明の効果】以上説明したように、本発明によれば、各入出力の用途に適合した用途別のバッファ管理を行うことが可能となり、システム全体の性能向上に寄与するところが多い。特に、例えば数十台～数千台というような多数の処理装置から構成される並列計算機システム等においては、各処理装置において限られたバッファ領域を効率よく管理し利用する技術が必要とされ、本発明によれば、このようなシステムにおいてバッファを利用した入出力処理の高速化と、バッファ領域の有効利用を

実現することが可能になる。

#### 【図面の簡単な説明】

【図1】本発明の構成例を示す図である。

【図2】本発明の実施例におけるファイルシステムの構成例説明図である。

【図3】本発明の実施例におけるバッファ領域の分割説明図である。

【図4】本発明の実施例におけるバッファブロックの管理情報説明図である。

10 【図5】本発明の実施例におけるスワップ制御説明図である。

【図6】本発明の実施例におけるスワップ中のブロックの状態遷移を示す図である。

【図7】本発明の実施例によるバッファ領域初期化処理フローチャートである。

【図8】本発明の実施例によるバッファブロック獲得処理フローチャートである。

【図9】本発明の実施例によるバッファブロック解放処理フローチャートである。

20 【図10】本発明の実施例によるスワップ処理フローチャートである。

#### 【符号の説明】

10 処理装置

11 ファイルシステム

12 ブロックサイズ指定手段

13 下限値指定手段

14 バッファ管理処理部

15 バッファ領域初期化処理部

16 バッファブロック獲得処理部

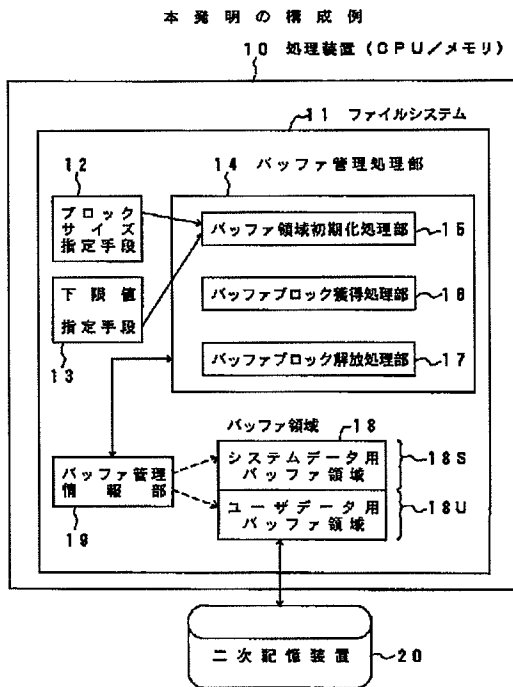
30 17 バッファブロック解放処理部

18 バッファ領域

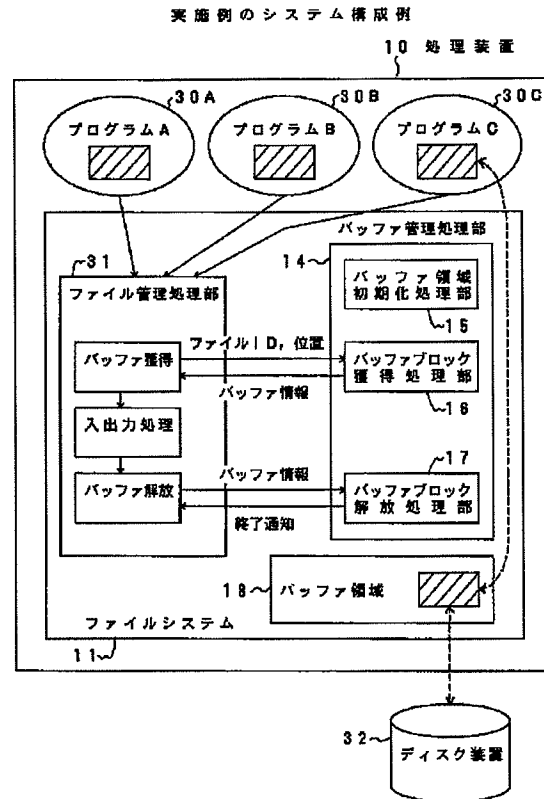
19 バッファ管理情報部

20 二次記憶装置

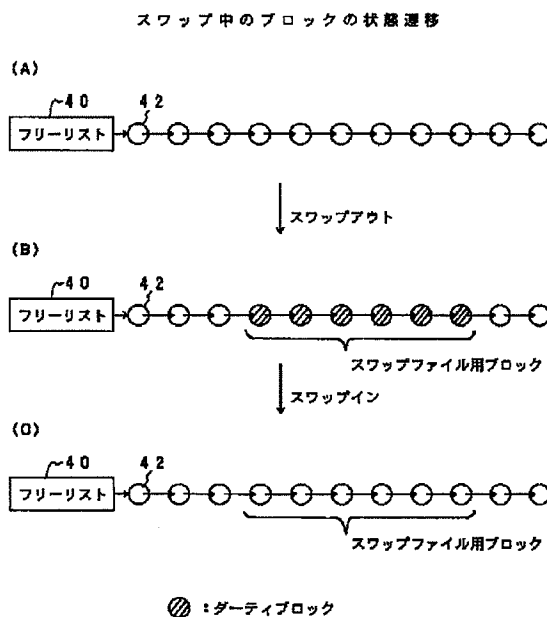
【図1】



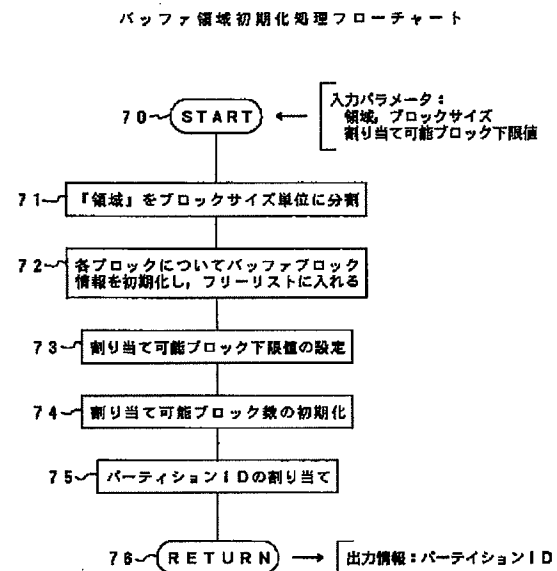
【図2】



【図6】

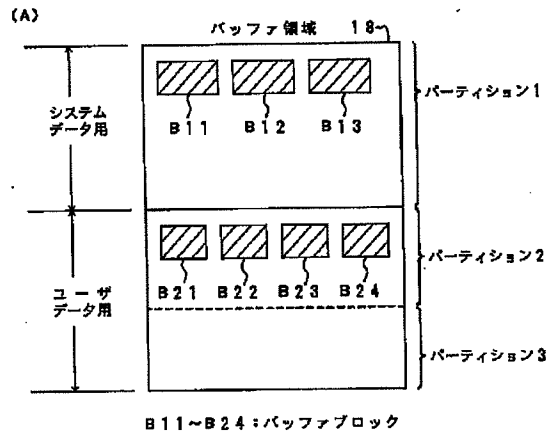


【図7】



【図3】

バッファ領域の分割説明図

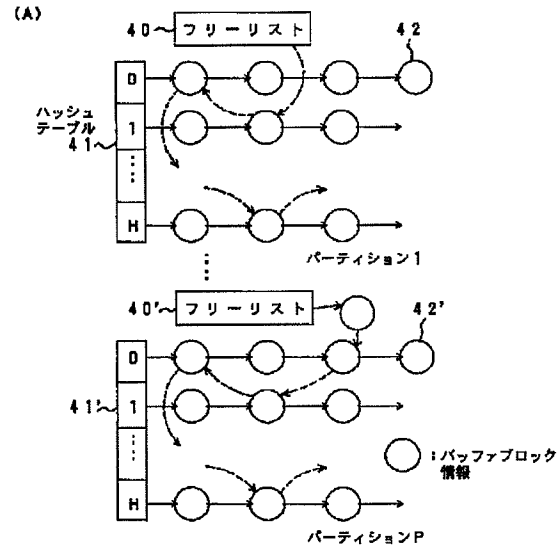


(B) パーティション管理情報

パーティション番号	ブロックサイズ	割り当て可能ブロック下限値
1	S1	L1
2	S2	L2
3	S3	L3

【図4】

バッファブロックの管理情報説明図



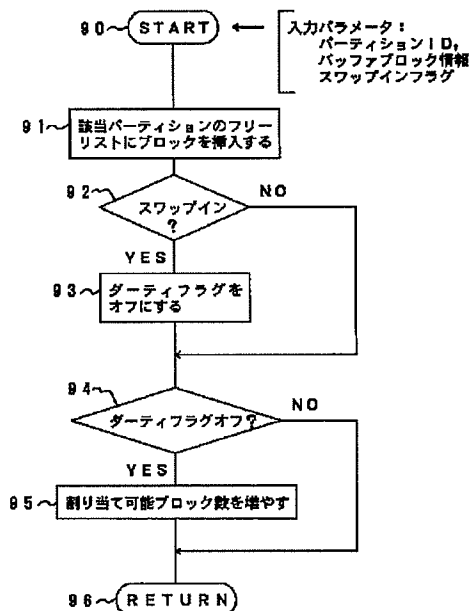
(B)

ファイルID	
ファイル内位置 (オフセット)	
バッファ内アドレス	
ダーティフラグ	データ有効フラグ
ハッシュキューのポインタ	
フリーリストのポインタ	
:	

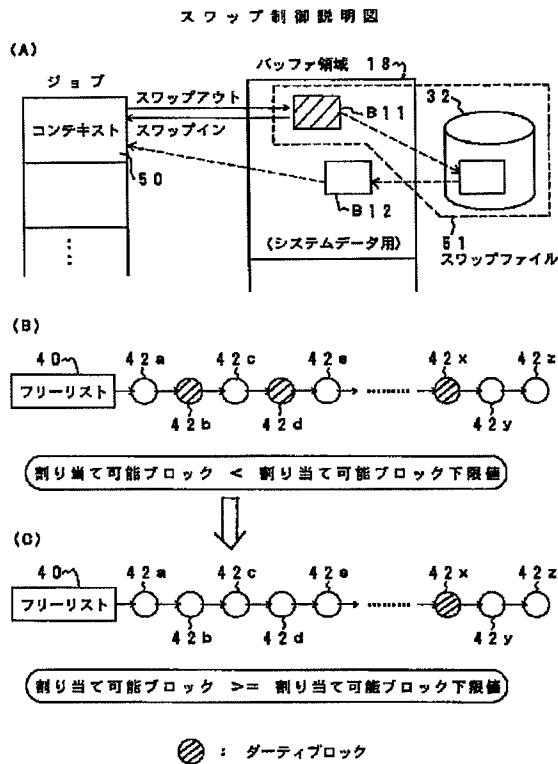
42 バッファブロック情報

【図9】

バッファブロック解放処理フローチャート

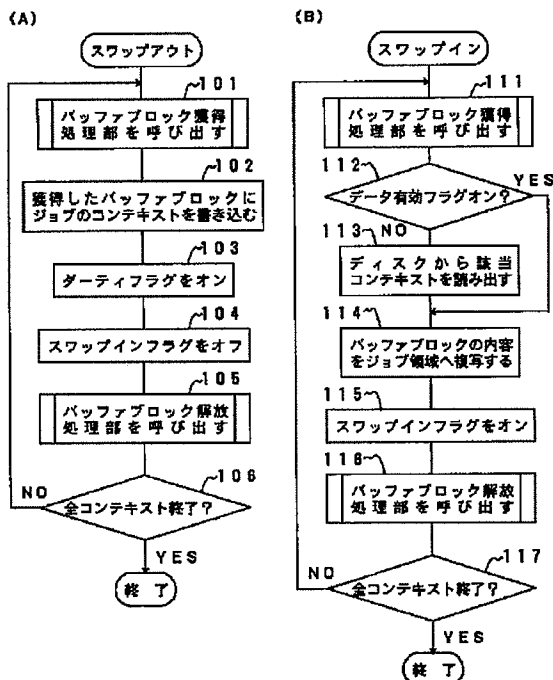


【図5】



【図10】

スワップ処理フローチャート



【図8】

